

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6285943号
(P6285943)

(45) 発行日 平成30年2月28日 (2018. 2. 28)

(24) 登録日 平成30年2月9日 (2018. 2. 9)

(51) Int. Cl.

F I

G 0 6 F 17/22 (2006. 01)

G 0 6 F 17/22 6 8 2

G 0 6 N 99/00 (2010. 01)

G 0 6 N 99/00 1 5 0

G 0 6 F 8/65 (2018. 01)

G 0 6 F 9/06 6 3 0 A

請求項の数 19 (全 20 頁)

(21) 出願番号 特願2015-537832 (P2015-537832)
 (86) (22) 出願日 平成25年10月17日 (2013. 10. 17)
 (65) 公表番号 特表2016-502701 (P2016-502701A)
 (43) 公表日 平成28年1月28日 (2016. 1. 28)
 (86) 国際出願番号 PCT/US2013/065497
 (87) 国際公開番号 W02014/062948
 (87) 国際公開日 平成26年4月24日 (2014. 4. 24)
 審査請求日 平成28年9月14日 (2016. 9. 14)
 (31) 優先権主張番号 13/653, 581
 (32) 優先日 平成24年10月17日 (2012. 10. 17)
 (33) 優先権主張国 米国 (US)

(73) 特許権者 314015767
 マイクロソフト テクノロジー ライセン
 シング, エルエルシー
 アメリカ合衆国 ワシントン州 9805
 2 レッドモンド ワン マイクロソフト
 ウェイ
 (74) 代理人 100079108
 弁理士 稲葉 良幸
 (74) 代理人 100109346
 弁理士 大貫 敏史
 (74) 代理人 100117189
 弁理士 江口 昭彦
 (74) 代理人 100134120
 弁理士 内藤 和彦

最終頁に続く

(54) 【発明の名称】 文字列変換の帰納的合成のための順位付け

(57) 【特許請求の範囲】

【請求項 1】

各々が部分表現を含むプログラム表現を含む候補変換プログラムを順位付けして、ユーザーによって入力される入力文字列から、1つまたは複数のユーザー供給される入力 - 出力例の各々と調和する、各々がユーザー所望の形式で出力文字列を生成する1つまたは複数の変換プログラムの順位付けされたグループを確立するために、コンピュータにより実装されるプロセスであって、

コンピュータを使用して、以下のプロセス動作、すなわち、

前記1つまたは複数のユーザー供給される入力 - 出力例から帰納的に合成された、各々が各ユーザー供給される入力例から各ユーザー供給される出力例によって提示される形式で出力文字列を生成する候補変換プログラムのセットを入力することと、

各候補変換プログラムについて、

より小さい方からより大きい方への順に、順位付け体系が確立された、前記候補変換プログラムの各部分表現について、前記部分表現についての前記順位付け体系を使用して尤度スコアを確立することと、

前記候補変換プログラムについて確立された前記部分表現の尤度スコアから、前記候補変換プログラムについての全体的な順位付けスコアを算出することとを実行させることを含む、コンピュータにより実装されるプロセス。

【請求項 2】

部分表現が、正規表現または位置表現またはアトミック表現または連結表現のいずれか

であり、前記順位付け体系が、前記正規表現、位置表現、アトミック表現、および連結表現の各々をこの順序で順位付けするために確立される、請求項 1 に記載のプロセス。

【請求項 3】

アトミック表現が、定数文字列表現または部分文字列表現のいずれかであり、アトミック表現についての順位付け体系を確立するプロセス動作が、

前記アトミック表現が定数文字列表現であるのか、または部分文字列表現であるのかを識別する動作と、

各タスクが複数の入力 - 出力例から成るタスクのトレーニングセットを入力する動作と、

前記トレーニングセットからトレーニングデータを生成する動作と、

各部分文字列表現および各定数文字列表現に尤度スコアを割り当てるために、機械学習分類子をトレーニングする動作を含む、請求項 2 に記載のプロセス。

【請求項 4】

各部分文字列表現に尤度スコアを割り当てるプロセス動作が、0 から 1 までの範囲における値を有する尤度スコアを割り当てることを含み、各定数文字列表現に尤度スコアを割り当てるプロセス動作が、所定の尤度スコア値を割り当てることを含み、請求項 3 に記載のプロセス。

【請求項 5】

前記タスクのトレーニングセットからトレーニングデータを生成するプロセス動作が、

前記トレーニングセット内の各タスクについて、

前記タスク内の入力 - 出力例ごとに、

前記入力 - 出力例から帰納的に合成され、各々が前記例内の入力文字列のタプルから前記例内の前記出力文字列を生成する、変換プログラムのセットを確立する動作と、

前記変換プログラムのセットから正のトレーニング部分表現のセットを確立する動作と、

前記変換プログラムのセットから負のトレーニング部分表現のセットを確立する動作と、

前記正および負のトレーニング部分表現のセット内の各部分文字列表現から前記トレーニングデータを生成する動作とを含む、請求項 3 に記載のプロセス。

【請求項 6】

前記正および負のトレーニング部分表現のセット内の各部分文字列表現からトレーニングデータを生成する前記プロセス動作が、

特徴の各々が、前記出力文字列内の部分文字列が前記正および負のトレーニング部分表現のセット内の部分文字列表現によって生成されるかを示す、1 つまたは複数の特徴のセットを入力する動作と、

前記正および負のトレーニング部分表現のセット内の各部分文字列表現について、各要素が、クラスラベルのために予約される要素を除いて、前記特徴のうちの 1 つにそれぞれ関連付けられる特徴値に対応する、要素の所定の順序を含む特徴ベクトルを生成する動作とを含み、前記特徴ベクトルの生成は、

前記要素の所定の順序に対応する順序における前記特徴の各々について、前記特徴が考慮中の前記部分文字列表現によって生成される前記出力文字列内の部分文字列において提示される場合、第 1 のバイナリ値を有する特徴値を、対応する特徴ベクトル要素に割り当て、前記特徴が考慮中の前記部分文字列表現によって生成される前記出力文字列中の部分文字列において提示されない場合、第 2 のバイナリ値を有する特徴値を、前記対応する特徴ベクトル要素に割り当てることと、

前記クラスラベルのために予約される前記特徴ベクトル要素について、前記部分文字列表現が前記正のトレーニング部分表現のセットに存在する場合、1 つのバイナリ値を割り当て、前記部分文字列表現が前記負のトレーニング部分表現のセットに存在する場合、異なるバイナリ値を割り当てることとを含む、請求項 5 に記載のプロセス。

【請求項 7】

前記 1 つまたは複数の特徴のセットは、

前記出力部分文字列の左端位置がトークンによって認識されること、
前記出力部分文字列の前記左端位置の文字が定数であること、
前記出力部分文字列の右端位置がトークンによって認識されること、
前記出力部分文字列の前記右端位置の文字が定数であること、
出力部分文字列はトークンを示すこと、
前記入力部分文字列の左端位置がトークンによって認識されること、
前記入力部分文字列の右端位置がトークンによって認識されること、
入力部分文字列はトークンであること、
前記出力部分文字列を取得するために実行されるケーシング、
前記部分文字列の長さ、
前記入力文字列の前記長さに対する前記部分文字列の相対的な長さ、
前記出力文字列の前記長さに対する前記部分文字列の相対的な長さ、または
構成要素となる位置表現の尤度スコア、
の少なくとも 1 つを含む、請求項 6 に記載のプロセス。

10

【請求項 8】

前記機械学習分類子は、サポートベクターマシン (SVM) バイナリ分類子である、請求項 6 に記載のプロセス。

【請求項 9】

タスク内の入力 - 出力例から確立される前記変換プログラムのセットから正のトレーニング部分表現のセットを確立する前記プロセス動作が、前記タスク内の前記入力 - 出力例から合成される変換プログラムの全てのセットの共通部分において見出される各部分表現を、正のトレーニング部分表現として指定することを含む、請求項 5 に記載のプロセス。

20

【請求項 10】

前記タスク内の前記入力 - 出力例から確立される前記変換プログラムのセットから負のトレーニング部分表現のセットを確立する前記プロセス動作が、正のトレーニング部分表現として識別されない前記変換プログラムのセット内の各部分表現を、負のトレーニング部分表現として指定することを含む、請求項 9 に記載のプロセス。

【請求項 11】

正規表現および位置表現について順位付け体系を確立する前記プロセス動作が、有限数の予め定義されたトークンを有するセットから値を取る出力部分文字列を生成する部分表現を、前記正のトレーニング部分表現のセット内のそれらの部分表現の頻度に基づいて、頻度に基づくスコアにマッピングする辞書を算出する動作を含む、請求項 5 に記載のプロセス。

30

【請求項 12】

連結表現について順位付け体系を確立する前記プロセス動作が、連結表現において見出される個々のアトミック表現についてのスコアを結合することを含む、請求項 2 に記載のプロセス。

【請求項 13】

連結表現において見出だされる個々のアトミック表現についてスコアを連結する前記プロセス動作が、前記スコアを乗算すること、または前記スコアを加算することの 1 つを含む、請求項 12 に記載のプロセス。

40

【請求項 14】

1 つまたは複数の変換プログラムの前記順位付けされたグループとして、所定の数の上位の候補変換プログラムを選択するプロセス動作をさらに含む、請求項 1 に記載のプロセス。

【請求項 15】

1 つまたは複数の変換プログラムの前記順位付けされたグループとして、所定の数の上位の候補変換プログラムを選択する前記プロセス動作が、前記最高位の候補変換プログラムを選択することを含む、請求項 14 に記載のプロセス。

50

【請求項 16】

各々のプログラムが部分表現を含むプログラム表現を含む候補変換プログラムを順位付けして、ユーザーによって入力される入力文字列から、1つまたは複数のユーザー供給される入力 - 出力例の各々と調和する、各々がユーザー所望の形式で出力文字列を生成する1つまたは複数の変換プログラムの順位付けされたグループを確立するためのシステムであって、

1つまたは複数のコンピューティング装置であって、前記コンピューティング装置は複数のコンピューティング装置がある場合はいつでもコンピューターネットワークを介して相互に通信する、1つまたは複数のコンピューティング装置と、

前記1つまたは複数のコンピューティング装置によって実行されるプログラムモジュールを有するコンピュータープログラムとを含み、前記1つまたは複数のコンピューティング装置は、

前記1つまたは複数のユーザー供給される入力 - 出力例から帰納的に合成された、各々が各ユーザー供給される入力例から各ユーザー供給される出力例によって提示される形式で出力文字列を生成する前記候補変換プログラムのセットを入力し、

各候補変換プログラムについて、

より小さい方からより大きい方への順に、順位付け体系が確立された、前記候補変換プログラムの各部分表現について、前記部分表現について確立された前記順位付け体系を使用して尤度スコアを確立し、

前記候補変換プログラムについて確立された前記部分表現の尤度スコアから、前記候補変換プログラムについての全体的な順位付けスコアを算出するように、前記コンピュータープログラムの前記プログラムモジュールにより誘導される、システム。

【請求項 17】

候補変換プログラムの順位付けの前に前記順位付け体系を確立するためのプログラムモジュールをさらに含み、各順位付け体系は前記部分表現が、ユーザーによって入力される入力文字列から、前記1つまたは複数のユーザー供給される入力 - 出力例の各々と調和する、ユーザー所望の形式で出力文字列を生成することができるプログラムの一部である尤度を示す尤度スコアを生成し、前記順位付け体系確立プログラムモジュールが、正規表現、位置表現、アトミック表現、および連結表現の各々をこの順序で順位付けする前記順位付け体系を確立するためのサブモジュールを含み、部分表現は、正規表現または位置表現またはアトミック表現または連結表現のいずれかである、請求項 16 に記載のシステム。

【請求項 18】

順位付け体系が確立された前記候補変換プログラムの各部分表現について、前記部分表現のために確立された前記順位付け体系を使用して尤度スコアを確立する前記プログラムモジュールは、

正規表現および位置表現を頻度に基づくスコアにマッピングする辞書を採用することを含む、正規表現および位置表現のための頻度に基づく順位付けを使用することと、

アトミック表現のための特徴に基づく順位付けを使用することであって、アトミック表現は定数文字列表現または部分文字列表現のどちらかであり、前記特徴に基づく順位付けは、

前記アトミック表現が定数文字列表現または部分文字列表現のどちらかあるか特定することと、

前記アトミック表現が定数文字列表現である場合はいつでも、所定の尤度スコア値を割り当てることと、

前記アトミック表現が部分文字列表現である場合にはいつでも、尤度スコアを各部分文字列表現に割り当てるようにトレーニングされた機械学習分類子を使用して尤度スコア値を割り当てることと、を含む、アトミック表現のための特徴に基づく順位付けを使用することと、

連結表現において見出される個々のアトミック表現についてのスコアを結合することを含む、連結表現を順位付けるためのパスに基づく順位付けを使用すること、のためのサ

10

20

30

40

50

ブモジュールを含む、請求項 17 に記載のシステム。

【請求項 19】

1 つまたは複数の変換プログラムの前記順位付けされたグループとして、所定の数の前記上位の候補変換プログラムを選択するためのプログラムモジュールをさらに含む、請求項 16 に記載のシステム。

【発明の詳細な説明】

【背景技術】

【0001】

[001] 世界中の数百万人もの人々が、データの記憶および操作のためにスプレッドシート等を使用する。これらのデータ操作シナリオは、大量の入力情報を 1 つの形式から別の形式へ転換して、所望の出力を生成することに関与することが多い。典型的には、これらのタスクは、手動で、または、エンドユーザーもしくはエンドユーザーのためのプログラムのいずれかによって作成される、小さい、しばしば 1 回限りのコンピュータプログラムを使用して達成される。

10

【0002】

[002] 別のアプローチは、コンピュータを採用して、所望のデータ変換を達成するためのプログラムを合成しようとする試みに関与した。プログラムを合成する 2 つの主要なアプローチ、すなわち、演繹法および帰納法が存在する。演繹的プログラム合成において、完全な上位の仕様は、翻訳の各工程が何らかの公理を使用して有効にされる、対応する下位のプログラムに翻訳される。このアプローチは、ユーザーが完全な仕様を提供することを必要とし、これは、幾つかの場合において、プログラム自体を書くことよりも潜在的に困難となり得る。これは、帰納的合成アプローチを、近年より普及させることとなった。帰納的プログラム合成アプローチにおいて、プログラムは、入力 - 出力例のセットから成る仕様などの不完全な仕様から合成される。帰納的プログラム合成アプローチは、下位のポインター操作コードからスプレッドシートのマクロに及ぶ、様々なドメインからのプログラムを合成するために近年使用されてきた。

20

【0003】

[003] 帰納的プログラム合成アプローチにおける仕様は、不完全であり、しばしば曖昧であるため、所与の仕様と調和する、基になるドメイン固有言語の異なるプログラムが存在する。曖昧さを除去し、所望のプログラムへ収束させるために、ユーザーは、付加的な入力 - 出力例を提供することによって、仕様を強化する必要がある。例の数は、ドメイン固有言語の表現度と正比例する。すなわち、言語が表現的であるほど、所望のプログラムに収束させるためにより多くの入力 - 出力例が必要とされる。

30

【0004】

[004] ドメイン固有言語は、ユーザーが望む大抵のタスクを表現するために表現的である必要があるが、同時に、所望のプログラムを学習するための面倒な数の入力 - 出力例をユーザーが提供することは期待できない。

【発明の概要】

【0005】

[005] 本明細書において説明される順位付け技法の実施形態は、一般に、帰納的プログラム合成手続を通じて生成された候補変換プログラムを、少数のユーザー供給された入力 - 出力例を使用して、順位付けすることに関与する。一実施形態において、コンピュータは、候補変換プログラムを順位付けして、1 つまたは複数の変換プログラムの順位付けされたグループを確立するために使用される。1 つまたは複数の変換プログラムの各々は、ユーザーによって入力された各入力文字列から、1 つまたは複数のユーザー供給された入力 - 出力例の各々に調和する、ユーザー所望の形式における出力文字列を生成する。

40

【0006】

[006] より詳細には、候補変換プログラムのセットが入力される。各変換プログラムは、ドメイン固有言語のプログラム表現から構成され、プログラム表現は、部分表現から構成されることに留意されたい。順位付け体系は、候補変換プログラムにおいて見出される

50

部分表現についてのオフライントレーニングデータから確立される。これらの順位付け体系の各々は、部分表現が、ユーザーによって入力された各入力文字列から、ユーザー供給される入力 - 出力例の各々と調和する、ユーザー所望の形式で出力文字列を生成することが可能なプログラムの一部である尤度を示す尤度スコアを生成する。順位付け体系が確立された各候補変換プログラムの各部分表現について、尤度スコアは、その部分表現について確立された順位付け体系を使用して確立される。尤度スコアは、より小さい部分表現からより大きい部分表現の順序で算出される。次いで、候補変換プログラムごとに、その候補変換プログラムについて確立された部分表現尤度スコアから、全体的な順位付けスコアが算出される。

【 0 0 0 7 】

10

[007] 一実施形態において、候補変換プログラムの各々に関連付けられる全体的な順位付けは、所定の数の上位の候補変換プログラムを選択するために使用される。次いで、これらの上位のプログラムは、1つまたは複数の変換プログラムの前述された順位付けされたグループとして指定される。1つの実装において、所定の数は1であり、従って、最高位の候補変換プログラムのみが選択および指定されることに留意されたい。

【 0 0 0 8 】

[008] 本概要は、詳細な説明においてさらに後述される概念の集まりを簡略化された形式で紹介するために提供されることに留意すべきである。本概要は、特許請求の範囲に記載された主題の重要な特徴または本質的な特徴を識別することを意図されたものではなく、特許請求の範囲に記載された主題の範囲を判定する補助として使用されることを意図されたものでもない。

20

【 0 0 0 9 】

[009] 本開示の具体的な特徴、態様、および利点は、以下の説明、添付の特許請求の範囲、および付属の図面に関して、より良く理解されるであろう。

【図面の簡単な説明】

【 0 0 1 0 】

【図1】[0010] 候補変換プログラムを順位付けするためのプロセスの一実施形態を一般的に概説するフロー図である。

【図2】[0011] アトミック表現のための順位付け体系を確立することに関与する、図1のプロセスの一部の実装を一般的に概説するフロー図である。

30

【図3】[0012] タスクのトレーニングセットからトレーニングデータを生成することに関与する、図2のプロセスの一部の実装を一般的に概説するフロー図である。

【図4】[0013] 正および負のトレーニング部分表現のセット内の各部分文字列からトレーニングデータを生成することに関与する、図3のプロセスの一部の実装を一般的に概説するフロー図である。

【図5】[0014] 順位付け体系が確立された候補変換プログラムの各部分表現についての尤度スコアを、その順位付け体系を使用して確立することに関与する、図1のプロセスの一部の実装を一般的に概説するフロー図である。

【図6】[0015] 最高位の候補変換プログラムを使用して、文字列変換を実行するためのプロセスの一実施形態を一般的に概説するフロー図である。

40

【図7】[0016] 人の名字を表す入力文字列を含む入力列と、文字列「Mr .」が名字の前に配置される対応する入力文字列名を含む1つの例示的な出力文字列を有する出力列とを有するスプレッドシートを描く表である。

【図8】[0017] 都市名を含む住所を含む入力列と、都市名に対応する、対応する入力文字列の一部を含む1つの例示的な出力文字列を有する出力列とを有するスプレッドシートを描く表である。

【図9】[0018] 本明細書において説明される順位付け技法の実施形態を実装するための例示的なシステムを構成する汎用コンピューティング装置を描く図である。

【発明を実施するための形態】

【 0 0 1 1 】

50

[0019] 順位付け技法の実施形態の以下の説明において、この説明の一部を形成し、本技法が実施され得る具体的な実施形態が例示として示される付属の図面への参照がなされる。本技法の範囲から逸脱することなく、他の実施形態が利用され、構造的な変更が行われ得ることが理解される。

【0012】

[0020] 明確さのために、本明細書において説明される順位付け技法の実施形態を説明する際には具体的な専門用語が行使され、これらの実施形態がそのように選ばれた具体的な用語に限定されることは意図されないことにも留意されたい。さらに、各具体的な用語は、同様の目的を実現するために概ね同様に動作する、その用語のあらゆる技術的な均等物を含むことが理解されるべきである。本明細書における、「一実施形態」、もしくは「別の実施形態」、もしくは「例示的な実施形態」、もしくは「代替的な実施形態」、または「1つの実装」、もしくは「別の実装」、もしくは「例示的な実装」、もしくは「代替的な実装」への言及は、その実施形態または実装に関連して説明される特定の特徴、特定の構造、または特定の特性が、本順位付け技法の少なくとも1つの実施形態に含まれ得ることを意味する。本明細書中の様々な箇所における「一実施形態において」、「別の実施形態において」、「例示的な実施形態において」、「代替的な実施形態において」、「1つの実装において」、「別の実装において」、「例示的な実装において」、「代替的な実装において」という表現の出現は、必ずしも全てが同じ実施形態または実装に言及しているわけではなく、別個または代替的な実施形態／実装は、他の実施形態／実装と相互に排他的でもない。またさらに、本順位付け技法の1つまたは複数の実施形態または実装を表すプロセスフローの順序は、いかなる特定の順序も本質的に示さず、本順位付け技法のいかなる限定も示唆しない。

【0013】

[0021] 本明細書において使用される「入力 - 出力例」という用語は、入力を形成する文字列のタプルおよび出力を形成する文字列に言及することにさらに留意されたい。入力は、ユーザーが変換することを希望する入力の例を表すのに対して、出力文字列は、ユーザーが入力から生成されることを希望する出力の例を表す。

【0014】

1.0 帰納的プログラム合成のための順位付け技法

[0022] 本明細書において説明される順位付け技法の実施形態は、統計的および機械的学習技法を使用して、帰納的プログラム合成において使用するための所望の順位付け関数を学習する。一般に、正および負の例のトレーニングデータセットは、入力 - 出力例の所与のセットから自動的に作成される。トレーニングデータセットから、ドメイン固有言語のプログラムにおける表現を尤度測度に割り当てる順位付け関数が学習される。次いで、これらの順位付け関数が使用されて、ごく少数の入力 - 出力例から、学習されたプログラムの尤度が算出される。

【0015】

1.1 文字列変換言語 L_S

[0023] このセクションでは、順位付け技法の実施形態を実装するための文字列変換言語が説明される。文字列変換言語 L_S についての構文は、以下の通りである。

```
Trace   expr e      := Concatenate ( f1 , . . . , fn )
Atomic  expr f      := SubStr ( vi , p1 , p2 )
                        | ConstStr ( s )
Position expr p     := Cpos ( k ) | Pos ( r1 , r2 , c )
Integer expr c      := k
Regular expr r      := TokenSeq ( T1 , . . . , Tm ) .
```

【0016】

[0024] トレース（または連結）表現 e は、アトミック表現 f_1, \dots, f_n の連結を示す。アトミック表現 f は、定数文字列表現 $\text{ConstStr}(s)$ または部分文字列表現 $\text{SubStr}(v_i, p_1, p_2)$ のいずれかを示し得る。定数文字列表現 Const

$\text{Str}(s)$ は、定数文字列「 s 」を生成する表現を示す。部分文字列表現 $\text{SubStr}(v_i, p_1, p_2)$ は、列 v_i に存在し、その左端位置および右端位置がそれぞれ位置ペア p_1 および p_2 によって表される入力文字列の部分文字列を表す。位置表現 $\text{CPos}(k)$ は、整数定数が正（または負）である場合、左側（または右側）からの所与の文字列中の k 番目のインデックスを指す。位置表現 $\text{Pos}(r_1, r_2, c)$ は、その左側および右側がそれぞれ正規表現 r_1 および r_2 に一致する入力文字列中の位置を示し、その文字列における正規表現の c 番目のそのような一致である。正式には、正規表現 r_1 は、 $s[0 \dots t-1]$ の何らかのサフィックスに一致し、 r_2 は、 $s[0 \dots l+1]$ の何らかのサフィックスに一致し、ただし、 $l = \text{Length}(s)$ であり、 t は、 c が正（または負）である場合、左側（または右側）から開始する c 番目のそのようなインデックス / 位置である。

10

【0017】

1. 2 L_s 表現の順位付けセット

[0025] L_s 表現の大きなセットを表すデータ構造は、以下の通りである。

【0018】

【数1】

$$\tilde{e} := \text{Dag}(\tilde{\eta}, \eta^s, \eta^t, \tilde{\xi}, W)$$

$$\tilde{f} := \text{SubStr}(v_i, \{p_j\}_j, \{p_k\}_k)$$

20

$$| \text{ConstStr}(s)$$

$$p := \text{CPos}(k) | \text{Pos}(r_1, r_2, c)$$

$$r := \text{TokenSeq}(T_1, \dots, T_m)$$

30

【0019】

[0026] データ構造は、多項式空間における表現の指数部つきの数値を表すために、ある表現の独立した部分表現を独立して維持する。例えば、部分文字列のセット $\text{SubStr}(v_i, \{p_j\}_j, \{p_k\}_k)$ は、左端位置の式 ($\{p_j\}_j$) および右端位置の式 ($\{p_k\}_k$) を独立して維持する。マッピング W は、 $\tilde{\eta}$ における各エッジを SubStr および ConstStr のセットにマッピングし、それによって、異なる長さの連結表現間での共有さえ実現する（開始ノード s から対象ノード t までの無閉路有向グラフにおける各パスは、そのパスにおける各エッジからの任意の表現を取り出して、それらを連結することによって取得された連結表現のセットを表す）。本明細書において説明される順位付け技法の実施形態は、この部分表現の独立性を維持し、部分表現のセットを独立して順位付けする。部分表現のセットを順位付けするための3つの技法、すなわち、（正規表現 r および位置表現 p を順位付けするための）頻度に基づく順位付け、（ SubStr 表現を順位付けするための）特徴に基づく順位付け、および（連結表現を順位付けするための）パスに基づく順位付けが使用される。

40

【0020】

1. 2. 1 頻度に基づく順位付け

[0027] 正規表現および位置表現（これらの双方は、何らかの有限のセットから値を取る）の場合、そのような表現を順位付けするために、頻度に基づく順位付けが実行される。

50

取り得る各表現値をトレーニングデータから推定されるその表現値の頻度にマッピングする辞書Dが存在する。正規表現の尤度スコアは、その頻度スコアに正比例するように定義される。位置表現の尤度スコアは、その頻度スコア、およびその構成要素である正規表現の尤度スコアの平方根の任意の線形関数となるように定義され得る。

【0021】

1.2.2 特徴に基づく順位付け

[0028] 無限のセットから値を取るアトミック表現fなどの表現の場合、特徴に基づく順位付けが実行される。特徴のセット (e) は、基となる入力 - 出力例から算出され、尤度は、機械学習手続から取得される。アトミック表現fについての定数文字列表現と部分文字列表現との間の選択を順位付けするために、サポートベクターマシン (SVM: support vector machine) バイナリ分類手続が、一実施形態において使用される。

10

【0022】

1.2.3 パスに基づく順位付け

[0029] アトミック表現から構成される連結表現の場合、パスに基づく順位付けが実行される。一実施形態において、このパスに基づく順位付けは、連結表現において見出される個々のアトミック表現についてのスコアを結合することに関与する。1つの実装において、アトミック表現についての尤度スコアは、これらを乗算することによって結合され、別の実装において、これらは互いに加算される。

【0023】

1.3 自動化されたトレーニングデータ生成

20

[0030] このセクションでは、頻度に基づく順位付け（または、それぞれ特徴に基づく順位付け）についての入力 - 出力例のセットから、正の例（または、正および負の例）から成るトレーニングデータを自動的に生成するための方法の一実施形態が説明される。各々が複数の入力 - 出力例 $\{(i_1, o_1), \dots, (i_n, o_n)\}$ から成るタスクの大きなセットが取得される。合成手続は、各入力 - 出力例 (i_k, o_k) と調和する全てのプログラム

【0024】

【数2】

$$\tilde{P}_k$$

30

のセットを学習する。次いで、正および負の表現値は、以下のように取得される。正の表現値は、セット

【0025】

【数3】

$$\tilde{P}_1 \cap \tilde{P}_2 \dots \cap \tilde{P}_n$$

に存在する表現値から構成され、負の表現値は、セット

【0026】

【数4】

$$\{\tilde{P}_k \setminus (\tilde{P}_1 \cap \tilde{P}_2 \dots \cap \tilde{P}_n) \mid 1 \leq k \leq n\}$$

40

に存在する表現値から構成される。

【0027】

1.3.1 正および負の表現値

[0031] 2つの無閉路有向グラフ (dags) D_k および D (D_1, D_2, \dots, D_n) が与えられると、課題は、2つのdagにおけるエッジを揃えて、正および負の表現値を算出することである。dagを揃えた後、これらの間の共通エッジは、正の表現値を構成し、 D_k には存在するが、 D には存在しないエッジは、負の表現値を構成する。

50

【 0 0 2 8 】

[0032] D A G プログラム D_k および D は、入力文字列 i_k に対して実行され、d a g ノードは、ラベル関数 $L: \text{int}$ を使用して、出力文字列 o_k のインデックスを用いて注釈付けされる。d a g の開始ノード s は、 $L(s) = 0$ となるように、インデックス 0 を用いて注釈付けされる。 $L(s_1) = 1$ であり、かつ、d a g エッジ上の表現 (s_1, s_2) が、入力文字列 i_k に対して実行される場合において、文字列 o_k を生成するとき、d a g におけるノード s_2 は、インデックス m (すなわち、 $L(s_2) = m$) を用いて注釈付けされる。双方の d a g のノードが注釈付けされると、同じラベルを有するノード間のエッジ上の表現は、比較のために収集される。d a g D_k において文字列 o_k [1 . . m] を生成する表現のセットは、

10

【 0 0 2 9 】

【数 5】

$$\tilde{e}_{l,m,k}$$

と示され、ただし、

【 0 0 3 0 】

【数 6】

$$\tilde{e}_{l,m,k} \equiv \bigcup_{\eta_1, \eta_2 \in \mathcal{D}_k} \tilde{e}(\eta_1, \eta_2), L(\eta_1) = l, L(\eta_2) = m$$

である。

20

【 0 0 3 1 】

【数 7】

$$\tilde{e}_{l,m,\wedge}$$

において出現する表現は、正の表現として示され、セット

【 0 0 3 2 】

【数 8】

$$\tilde{e}_{l,m,k} \setminus \tilde{e}_{l,m,\wedge}$$

において出現する表現は、負の表現として示される。正および負の表現値のセットは、各入力 - 出力例のペアについて取得される。

30

【 0 0 3 3 】

1 . 3 . 2 頻度に基づくトレーニングデータ生成

[0033] 正規表現および位置表現を順位付けする場合、これらの双方は有限のセットから値を取るため、頻度に基づく順位付けアプローチが実行される。トークンシーケンス表現についてのトークンシーケンスとこれらの頻度スコアのデータベースも、作成される。頻度は、正のトレーニング表現のセットから推定される。頻度は、ある表現が発生する各異なるコンテキストについても推定され得る。例えば、ある位置表現は、左端位置表現または右端位置表現のどちらかであり得る(その位置表現が S u b S t r 表現の第 1 の引数として発生するか、または第 2 の引数として発生するかを意味する)。正規表現も、ある位置表現内の 2 つの異なるコンテキストにおいて発生する。

40

【 0 0 3 4 】

1 . 3 . 3 特徴に基づくトレーニングデータ生成

[0034] 文字列変換を学習する際の主要な曖昧さのうちの 1 つは、出力文字列中の部分文字列が定数文字列であるのか、または入力文字列の部分文字列であるのかに関する決定を行うことに由来する。そのような決定を行うことは、位置ペアおよび定数文字列についての取り得る値だけでなく、入力文字列および出力文字列にも依存するため、無限に多くの取り得る入力文字列および出力文字列が存在することから、この場合において頻度に基づく順位付けを使用することは不可能である。代わりに、一実施形態において、特徴に基づく順位付けアプローチが使用されて、S u b S t r 表現と C o n s t S t r 表現との間で

50

選択がされる。より詳細には、特徴のセットが、各 `SubStr` 表現について定義される。使用され得る特徴の例は、以下を含むが、以下に限定されない。

- a) `IsOutputLeftTok` : 出力部分文字列の左端位置がトークンによって認識され得るかを示すブール値
- b) `IsOutputLeftConstant` : 出力部分文字列の左端位置の文字が定数であるかを示すブール値
- c) `IsOutputRightTok` : 出力部分文字列の右端位置がトークンによって認識され得るかを示すブール値
- d) `IsOutputRightConstant` : 出力部分文字列の右端位置の文字が定数であるかを示すブール値
- e) `IsOutputTok` : 出力部分文字列がトークンを意味するかを意味するブール値
- f) `IsInputLeftTok` : 入力部分文字列の左端位置がトークンによって認識され得るかを意味するブール値
- g) `IsInputRightTok` : 入力部分文字列の右端位置がトークンによって認識され得るかを意味するブール値
- h) `IsInputTok` : 入力部分文字列がトークンであることを意味するブール値
- i) `Casing` : 出力部分文字列を取得するために実行されるケーシング
- j) `LenSubstring` : 部分文字列の長さ
- k) `RelLenInSubstring` : 入力文字列の長さに対する部分文字列の相対的な長さ (`lenSubstring/lenInputString`)
- l) `RelLenOutSubstring` : 出力文字列の長さに対する部分文字列の相対的な長さ (`lenSubstring/lenOutputString`)
- m) 頻度に基づく順位付け体系を使用して推定される、構成要素となる位置表現の尤度スコア

【0035】

[0035] これらの特徴全ては、一定の $O(1)$ 時間で算出され得る。各正および負の `SubStr` 表現について、特徴ベクトルは、クラスラベル（例えば、正の表現について +1、および負の表現について 0）と共に算出される。次いで、既製のサポートベクターマシン（SVM）手続が、バイナリ分類子を学習するために使用され得る。より詳細には、作成された各特徴ベクトルは、各要素（クラスラベルのために予約される要素を除く）が前述された特徴のうちの 1 つにそれぞれ関連付けられる特徴値に対応する所定の順序を有する。要素の所定の順序に対応する順序における前述された特徴の各々について、その特徴が考慮中の `SubStr` 表現において提示される場合、第 1 のバイナリ値（例えば、1）を有する特徴値が、対応する特徴ベクトル要素に割り当てられる。その特徴が提示されない場合、第 2 のバイナリ値（例えば、0）を有する特徴値が、対応する特徴ベクトル要素に割り当てられる。また、前述されたクラスラベルのために予約される特徴ベクトル要素は、`SubStr` 表現が正の `SubStr` 表現である場合、1 つのバイナリ値（例えば、1）を割り当てられ、`SubStr` 表現が負の `SubStr` 表現である場合、別のバイナリ値（例えば、0）を割り当てられる。

【0036】

1.4 順位付けプログラム

[0036] このセクションでは、`dag` によって表されるプログラムがどのように順位付けされるかが説明されるであろう。

【0037】

1.4.1 Dag エッジ表現の順位付け

[0037] `dag` の各エッジは、`SubStr` 表現および `ConstStr` 表現のセットから成る。一実施形態において、特徴に基づく順位付けは、+1 と 0 との間の尤度スコアを各 `SubStr` 表現に割り当てる一方で、`ConstStr` 表現の尤度スコアは、0.5 になるように決められる。

【 0 0 3 8 】

1 . 4 . 2 D a g パスの順位付け

[0038] d a g Dにおける各パスは、所与の出力例に適合する幾つかのプログラムを表す。エッジeの尤度スコア $w(e)$ は、そのエッジにおける任意のS u b S t r表現またはC o n s t S t r表現の最も高い尤度スコアとして定義される。パスの尤度スコアは、そのパス上の対応するエッジ表現の尤度スコアを乗算する（加算する）ことによって算出されるように、 $w(p) = \prod_{e \in edges(p)} w(e)$ として定義される。次いで、ダイクストラの最短パス手続が使用されて、d a gにおける最高位のパスが算出され、最高位の出力文字列を生成するために、このパスが実行される。パスの尤度スコアは、全てのエッジの尤度スコアの積／和、エッジの数、任意のエッジの最小／最大の尤度スコアなど、様々な特徴のより洗練された関数とすることもできる。さらに、この関数は、機械的学習技法を使用することによって学習されることもできる。

10

【 0 0 3 9 】

1 . 5 候補変換プログラムを順位付けするための例示的なプロセス

[0039] 本明細書において説明される順位付け技法の実施形態の前述の態様は、1つの一般的な実装において、図1において概説されるプロセスによって実現され得る。より詳細には、（後述の例示的な動作環境において説明されるコンピューティング装置のいずれかなどの）コンピューターは、候補変換プログラムを順位付けして、1つまたは複数の変換プログラムの順位付けされたグループを確立するために使用される。1つまたは複数の変換プログラムの各々は、ユーザーによって入力された各入力文字列から、1つまたは複数のユーザー供給された入力 - 出力例の各々と調和する、ユーザー所望の形式における出力文字列を生成する。より詳細には、候補変換プログラムのセットは、入力として受信される（プロセス動作100）。候補変換プログラムの各々は、1つまたは複数のユーザー供給された入力 - 出力例から、従来の方法を使用して帰納的に合成され、各ユーザー供給された入力例から各ユーザー供給された出力例によって提示される形式で出力文字列を生成する。既述したように、各変換プログラムは、部分表現から構成されるプログラム表現から構成されることに留意されたい。次に、（より小さい部分表現からより大きい部分表現まで）順位付け体系が確立された各候補変換プログラムの各部分表現について、尤度スコアは、その部分表現について確立された順位付け体系を使用して確立される（プロセス動作102）。順位付け体系の各々は、部分表現が、ユーザーによって入力された各入力から、ユーザー供給された入力 - 出力例の各々と調和する、ユーザー所望の形式で出力文字列を生成することが可能なプログラムの一部である尤度を示す尤度スコアを生成することに留意されたい。次いで、全体的な順位付けスコアが、候補変換プログラムごとに、その候補変換プログラムについて確立された部分表現尤度スコアから算出される（プロセス動作104）。

20

30

【 0 0 4 0 】

[0040] 色々な種類の順位付け体系が、色々な種類の表現について使用される。既述したように、部分表現は、正規表現または位置表現またはアトミック表現または連結表現のいずれかであり得る。一実施形態において、正規表現および位置表現は、頻度に基づく順位付けを採用する（このタイプの順位付けは、トークンの有限なセットからの値を取る表現に適用可能であるため）。また、特徴に基づく順位付けは、アトミック表現について使用され、パスに基づく順位付けは、連結について使用される。効率のために、一実施形態において、順位付け体系は、まず正規表現に、次いで位置表現に、次いでアトミック表現に、最後に連結表現について確立されることにも留意されたい。

40

【 0 0 4 1 】

[0041] アトミック表現の場合において、既に示されたように、これらは、定数文字列表現または部分文字列表現のどちらかであり得る。これを考慮して、図2を参照すると、一実施形態においてアトミック表現についての順位付け体系を確立することは、まず、アトミック表現が定数文字列表現であるのか、または部分文字列表現であるのかを識別することに関与する（プロセス動作200）。次いで、タスクのトレーニングセットが入力され

50

る（プロセス動作202）。各タスクセットは、複数の入力 - 出力例から構成される。トレーニングデータは、トレーニングセットから生成され（プロセス動作204）、尤度スコアを各部分文字列表現および各定数文字列表現に割り当てるべく、機械学習分類子をトレーニングするために使用される（プロセス動作206）。一実施形態において、機械学習分類子は、サポートベクターマシン（SVM）バイナリ分類子であり、部分表現が意図されるプログラムの一部である可能性がどのくらいあるかに応じて0から1までの範囲における値を有する各部分文字列表現に尤度スコアを割り当てるためにトレーニングされる。また、分類子は、所定の尤度スコア値（例えば、0.5）を各定数文字列表現に割り当てる。

【0042】

[0042] タスクのトレーニングセットからトレーニングデータを生成することに関して、1つの実装において、これは、トレーニングセット内の各タスクについて以下に關与する。図3を参照すると、タスクにおいて、以前に選択されていない入力 - 出力例が選択される（プロセス動作300）。変換プログラムのセットは、選択された入力 - 出力例から帰納的に合成される（プロセス動作302）。変換プログラムの各々は、その例の入力における入力文字列のタプルから、その例の出力において提示される形式で出力文字列を生成する。入力文字列のタプルとは、入力 - 出力例の出力を生成するために使用される入力 - 出力例における入力の1つまたは複数の部分をいう。次に、正のトレーニング部分表現のセットが、変換プログラムのセットから確立され（プロセス動作304）、負のトレーニング部分表現のセットも、変換プログラムのセットから確立される（プロセス動作306）。一実施形態において、変換プログラムのセットから正のトレーニング部分表現を確立することは、タスク内の入力 - 出力例から合成される変換プログラムの全てのセットの共通部分において見出される各部分表現を、正のトレーニング部分表現に指定することに関与する。一方で、変換プログラムのセットから負のトレーニング部分表現を確立することは、正のトレーニング部分表現として識別されない変換プログラムのセット内の各部分表現を、負のトレーニング部分表現として指定することに関与する。

【0043】

[0043] この時点で、正規表現および位置表現についての順位付け体系が確立される。より詳細には、有限な数の予め定義されたトークンを有するセットから値を取る出力部分文字列を生成する部分表現（すなわち、正規表現および位置表現）を、正のトレーニング部分表現のセット内のそれらの部分表現の頻度に基づいて、頻度に基づくスコアにマッピングする辞書が生成される（プロセス動作308）。

【0044】

[0044] 次に、部分表現の正および負のトレーニングセットの双方からの各部分文字列表現について、トレーニングデータが、特徴に基づく順位付けのために生成される（プロセス動作310）。次いで、タスク内の全ての入力 - 出力例が選択および処理されたかが判定される（プロセス動作312）。タスク内の全ての入力 - 出力例が選択および処理されていない場合、プロセス動作300～312が繰り返される。そうでない場合、本手続は、選択されたタスクについて終了する。前述の手続は、残りのタスク全ておよびそれらの入力 - 出力例について繰り返される。

【0045】

[0045] 特徴に基づく順位付けのために、正および負のトレーニング部分表現のセット内の各部分文字列表現からトレーニングデータを生成する前述のプロセス動作に関して、これは、図4に示されるような一実施形態において達成される。まず、（既述された特徴などの）1つまたは複数の特徴のセットが入力される（プロセス動作400）。特徴の各々は、出力文字列中の部分文字列が正および負のトレーニング部分表現のセット内の部分文字列表現のうちの1つによって生成されるかを示す。次いで、正および負のトレーニング部分表現のセット内の各部分文字列表現について、特徴ベクトルが生成される（プロセス動作402）。既述したように、各特徴ベクトルは、クラスラベルのために予約された要素を除いて、各要素が前述された特徴のうちの1つにそれぞれ関連付けられる特徴値に対

10

20

30

40

50

応する、要素の所定の順序である。一実施形態において、要素の所定の順序に対応する順序における特徴の各々について、その特徴が考慮中の部分文字列表現によって生成される出力文字列中の部分文字列において提示される場合、各特徴ベクトルは、第1のバイナリ値（例えば、1）を有する特徴値を割り当てることによって生成される。より具体的には、特徴値は、提示される特徴に関連付けられる特徴ベクトル要素に割り当てられる。また、第2のバイナリ値（例えば、0）を有する特徴値は、考慮中の特徴が部分文字列表現によって生成される出力文字列中の部分文字列において提示されない場合、その特徴に関連付けられる特徴ベクトル要素に割り当てられる。さらに、前述されたクラスラベルのために予約される特徴ベクトル要素は、部分文字列表現が正のトレーニング部分表現のセット内にある場合、1つのバイナリ値（例えば、1）を割り当て、部分文字列表現が負のトレーニング部分表現のセット内にある場合、異なるバイナリ値（例えば、0）を割り当てる。

10

【0046】

[0046] 連結表現について順位付け体系を確立することに関して、一実施形態において、これは、連結表現において見出される個々のアトミック表現についての尤度スコアを結合することに関与する。1つの実装において、この結合は、スコアを乗算することによって達成され、かつ1つの実装では、この結合は、スコアを加算することによって達成される。

【0047】

[0047] 前記を考慮すると、一実施形態において、順位付け体系を使用してその順位付け体系が確立された候補変換プログラムの各部分表現についての尤度スコアを確立することは、以下に関与する。図5を参照すると、まず、各正規表現についての尤度が、頻度に基づく順位付け体系を使用して識別される（プロセス動作500）。次いで、各位置表現についての尤度も、頻度に基づく順位付け体系を使用して識別される（プロセス動作502）。既述したように、この頻度に基づく順位付けは、正規表現および位置表現を頻度に基づくスコアにマッピングする辞書を採用することに関与し得る。次に、各アトミック表現の尤度が、特徴に基づく順位付け体系を使用して識別される（プロセス動作504）。既述したように、アトミック表現は、定数文字列表現または部分文字列表現のどちらかである。アトミック表現が定数文字列表現である場合、特徴に基づく順位付けは、所定の尤度スコア値（例えば、0.5）を割り当てることに関与し得る。また、アトミック表現が部分文字列表現である場合、特徴に基づく順位付けは、尤度スコアを各部分文字列表現に割り当てるようにトレーニングされた機械学習分類子を使用して、（例えば、0から1までに及ぶ）尤度スコア値を割り当てることに関与し得る。次いで、各連結表現の尤度は、パスに基づく順位付け体系を使用して識別される（プロセス動作506）。既述したように、このパスに基づく順位付けは、連結表現において見出される個々のアトミック表現についてのスコアを結合することに関与し得る。

20

30

【0048】

[0048] 候補変換プログラムの各々について確立される、前述された全体的な順位付けに関して、これらの順位付けを多様な手法で使用することが可能である。例えば、一実施形態において、候補変換プログラムの各々に関連付けられる順位付けは、所定の数の上位の候補変換プログラムを選択し、これら上位のプログラムを1つまたは複数の変換プログラムの前述された順位付けされたグループとして指定するために使用される。1つの実装において、所定の数は1であり、従って、最高位の候補変換プログラムのみが選択および指定されることに留意されたい。

40

【0049】

1.6 最高位の変換プログラムを使用して文字列変換を実行するための例示的なプロセス

[0049] 最高位の候補変換プログラムを選択することに関与する後者の実施形態に関して、これは、コンピューティング装置が文字列変換を実行する適用において使用され得る。より詳細には、図6を参照すると、1つの一般的な実装において、1つまたは複数のユー

50

ザー供給された入力 - 出力例から帰納的に合成された候補変換プログラムのセットが入力される（プロセス動作 600）。既述したように、候補変換プログラムの各々は、各ユーザー供給された入力例から各ユーザー供給された出力例によって提示される形式で出力文字列を生成する。

【0050】

[0050] 入力されると、候補変換プログラムのセットは、順位付けされて、最高位の変換プログラムが識別される（プロセス動作 602）。次いで、この最高位の変換プログラムが、ユーザー供給された入力文字列に適用されて、出力文字列が生成される（プロセス動作 604）。

【0051】

1.7 例示的な変換シナリオ

[0051] このセクションでは、2つの例示的な変換シナリオが提供される。本明細書において説明される順位付け技法の実施形態は、これらのシナリオにおいて見出される少数の入力 - 出力例に基づいて、入力される各入力文字列から所望の出力文字列を生成することが可能な変換プログラムを確立するために採用され得る。

【0052】

[0052] 第1の例示的なシナリオにおいて、スプレッドシートのユーザーは、入力列において一連の名前を有し、各名前の前に敬称「Mr.」を追加することを希望する。ユーザーは、図7に示される入力 - 出力例を提供した。

【0053】

[0053] 課題は、出力文字列「Mr. Roger」中のどの部分文字列が定数文字列であり、どれが入力文字列「Roger」の部分文字列であるのかを決定することによって、所与の入力 - 出力例から所望の変換プログラムを学習することである。この場合には、出力部分文字列 $o_1[0..0]$ Mが入力文字列中に存在しないため、定数文字列でなければならないことが推論され得る。しかし、出力部分文字列 $o_1[1..1]$ rは、入力文字列中の2つの異なる部分文字列（ $i_1[0..0]$ Rおよび $i_1[4..4]$ r）に由来し得る。本明細書において説明される順位付け技法の実施形態は、全ての取り得る表現を学習して、（i）部分文字列 $i_1[0..0]$ を抽出し、「r」から「R」を生成するための小文字演算を実行するための位置ペア表現と、（ii）部分文字列 $i_1[4..4]$ および定数文字列表現「r」を抽出するための位置ペア表現とを含む出力文字列中の部分文字列「r」を算出する。入力文字列に存在しない出力文字列の部分文字列は、定数文字列であることが保証されるのに対して、入力文字列に存在しない部分文字列は、定数文字列または入力文字列の部分文字列のどちらかであり得る（経験は、これらが入力文字列に由来する可能性が高いことを示しているが）。例えば、出力部分文字列 $o_1[4..8]$ Rogerは、定数文字列「Roger」よりも、入力文字列 i_1 に由来する可能性が高い。同様の議論を使用すると、出力文字列中の部分文字列 $o_1[1..1]$ rも、入力文字列中の2つの位置に存在し、そこに由来する可能性が高い。しかし、この例において、プログラムの所望の振る舞いは定数文字列「Mr.」を各入力文字列の前に追加することであるため、「r」は定数文字列である必要がある。1つの入力 - 出力例から所望の変換を学習するために、本明細書において説明される順位付け技法の実施形態は、定数文字列表現を出力部分文字列「r」についての位置ペア表現よりも上位に順位付けする必要がある。「r」を定数文字列として順位付けする際に役立つ特徴のうちの幾つかは、以下を含む。

a) 部分文字列の長さ：部分文字列「r」の長さは1であるため、部分文字列「r」が入力部分文字列である可能性は低い。

b) 部分文字列の相対的な長さ：出力文字列と比較した部分文字列「r」の相対的な長さも非常に小さく、すなわち、 $1/10$ である。

c) 定数の隣接文字：「r」の隣接文字「M」および「.」は、双方ともに定数表現である。

【0054】

[0054] 第2の例示的なシナリオにおいて、スプレッドシートのユーザーは、列中に一連の住所を有し、この住所から都市名を抽出することを希望する。ユーザーは、図8に示される入力 - 出力例を提供した。

【0055】

[0055] この場合には、本明細書において説明される順位付け技法の実施形態は、入力文字列「243 Flyer Drive, Cambridge, MA 02145」から部分文字列「Cambridge」を抽出するために、100を超える異なるSubStr式を学習し得る。そのうちの幾つかは、以下の通りである。

- a) p_1 : 3番目のアルファベットトークン文字列を抽出する。
- b) p_2 : 4番目の英数字のトークン文字列を抽出する。
- c) p_3 : コンマおよび空白のトークンから成る1番目と2番目のトークンシーケンスの間の部分文字列を抽出する。
- d) p_4 : 3番目の大文字トークン(包括的)と左から2番目のコンマトークンとの間の部分文字列を抽出する。

【0056】

[0056] 位置ペア式 p_1 および p_2 を学習することに関する問題は、入力文字列「512 Wright Ave, Los Angeles, CA 78911」について、所望の出力ではない出力文字列「Los」を生成することである。また、位置ペア式 p_4 は、入力文字列「64 128th St, Seattle, WA 98102」から、いかなる出力文字列も生成しない。他方、位置ペア式 p_3 は、位置ペア式の各々について所望の出力文字列(すなわち、それぞれCambridge、Los Angeles、SeattleおよびSan Francisco)を生成する。そのため、本明細書において説明される順位付け技法の実施形態は、位置ペア式 p_3 をその他の位置ペア式よりも上位に順位付けして、1つの入力 - 出力例から所望の出力文字列を生成する。

【0057】

2.0 例示的な動作環境

[0057] 本明細書において説明される順位付け技法の実施形態は、多数のタイプの汎用または専用コンピューティングシステム環境または構成内で動作可能である。図9は、本明細書において説明されるような様々な実施形態および順位付け技法の実施形態の要素が実装され得る汎用コンピューターシステムの簡略化された例を例示する。図9において破線または点線によって表されるいかなる四角形も、簡略化されたコンピューティング装置の代替的な実施形態を表すこと、および、以下に説明されるように、これらの代替的な実施形態のうちのいずれかまたは全ては、本文書全体にわたって説明される他の代替的な実施形態と組み合わせて使用され得ることが留意されるべきである。

【0058】

[0058] 例えば、図9は、簡略化されたコンピューティング装置10を示す、一般的なシステム図を示す。そのようなコンピューティング装置は、典型的には、パーソナルコンピューター、サーバーコンピューター、ハンドヘルドコンピューティング装置、ラップトップまたはモバイルコンピューター、携帯電話およびPDAなどの通信装置、マルチプロセッサシステム、マイクロプロセッサベースのシステム、セットトップボックス、プログラム可能な家庭用電化製品、ネットワークPC、ミニコンピューター、メインフレームコンピューター、オーディオまたはビデオメディアプレーヤー等を含むが、これらに限定されない、少なくとも何らかの最低限の計算能力を有する装置において見出され得る。

【0059】

[0059] 本明細書において説明される順位付け技法の実施形態を装置に実装させるために、装置は、基本的な計算演算を可能にするための十分な計算能力およびシステムメモリを有するべきである。特に、図9によって例示されるように、計算能力は、一般に、1つまたは複数の処理装置12によって例示され、1つまたは複数のGPU14も含み得、これらの一方または双方がシステムメモリ16と通信する。一般的なコンピューティング装置の処理装置12は、DSP、VLIW、もしくは他のマイクロコントローラーなどの特殊

10

20

30

40

50

マイクロプロセッサであってもよく、または、マルチコアCPUにおける特殊なGPUベースのコアを含む、1つまたは複数の処理スコアを有する従来のCPUとし得ることに留意されたい。

【0060】

[0060] また、図9の簡略化されたコンピューティング装置は、例えば、通信インターフェイス18などの他のコンポーネントも含み得る。図9の簡略化されたコンピューティング装置は、1つまたは複数の従来のコンピューター入力装置20（例えば、ポインティング装置、キーボード、オーディオ入力装置、ビデオ入力装置、触覚入力装置、有線または無線データ送信を受信するための装置等）も含み得る。図9の簡略化されたコンピューティング装置は、例えば、1つまたは複数の従来のディスプレイ装置24および他のコンピューター出力装置22（例えば、オーディオ出力装置、ビデオ出力装置、有線または無線データ送信を送信するための装置等）などの、他の任意選択的なコンポーネントも含み得る。汎用コンピューターのための典型的な通信インターフェイス18、入力装置20、出力装置22、および記憶装置26は、当業者には周知であり、本明細書において詳細に説明されないことに留意されたい。

10

【0061】

[0061] 図9の簡略化されたコンピューティング装置は、多様なコンピューター読取可能な媒体も含み得る。コンピューター読取可能な媒体は、コンピューター10によって記憶装置26を介してアクセスされ得る任意の利用可能な媒体とすることができ、コンピューター読取可能な命令もしくはコンピューター実行可能な命令、データ構造、プログラムモジュール、または他のデータなどの情報の記憶のために、取り外し可能28および/または取り外し不可能30である揮発性媒体および不揮発性媒体の双方を含む。コンピューター読取可能な媒体は、コンピューター記憶媒体および通信媒体を含み得る。コンピューター記憶媒体は、DVD、CD、フロッピーディスク、テープドライブ、ハードディスクドライブ、光学式ドライブ、ソリッドステートメモリドライブ、RAM、ROM、EEPROM、フラッシュメモリもしくは他のメモリ技術、磁気カセット、磁気テープ、磁気ディスク記憶装置もしくは他の磁気記憶装置、または所望の情報を記憶するために使用されることができ、かつ、1つまたは複数のコンピューティング装置によってアクセスされることができる任意の他の装置などの、有形のコンピューター読取可能なまたは機械読取可能な媒体または記憶装置をいう。

20

30

【0062】

[0062] コンピューター読取可能なまたはコンピューター実行可能な命令、データ構造、プログラムモジュール等などの情報の保持は、1つもしくは複数の変調データ信号もしくは搬送波を符号化するための多様な前述の通信媒体のいずれか、または、他の伝送機構もしくは通信プロトコルを使用することによっても達成され得、任意の有線もしくは無線の情報配信機構を含む。「変調データ信号」または「搬送波」という用語は、一般的に、その特性のうちの1つまたは複数が情報を信号に符号化するような手法で設定または変更された信号をいうことに留意されたい。例えば、通信媒体は、1つまたは複数の変調データ信号を搬送する有線ネットワークまたは直接有線接続などの有線媒体、ならびに音響、RF、赤外線、レーザー、および1つまたは複数の変調データ信号または搬送波を送信および/または受信するための他の無線媒体などの無線媒体を含む。上記のうちの任意のものの組み合わせも、通信媒体の範囲内に含まれるべきである。

40

【0063】

[0063] さらに、本明細書において説明される様々な順位付け技法の実施形態のうちの一部または全部を具現化するソフトウェア、プログラムおよび/もしくはコンピュータープログラム製品、またはこれらの一部は、コンピューターまたは機械読取可能な媒体または記憶装置および通信媒体の任意の所望の組み合わせから、コンピューター実行可能な命令または他のデータ構造の形式において、記憶され、受信され、送信され、または読み出され得る。

【0064】

50

[0064] 最後に、本明細書において説明される順位付け技法の実施形態は、コンピューティング装置によって実行される、プログラムモジュールなどのコンピューター実行可能な命令の一般的なコンテキストにおいてさらに説明され得る。一般に、プログラムモジュールは、特定のタスクを実行し、または特定の抽象データ型を実装するルーチン、プログラム、オブジェクト、コンポーネント、データ構造等を含む。本明細書において説明される実施形態は、1つもしくは複数の通信ネットワークを通じてリンクされる、タスクが1つもしくは複数の遠隔処理装置によって実行される分散コンピューティング環境において、または1つもしくは複数の装置のクラウド内でも実施され得る。分散コンピューティング環境において、プログラムモジュールは、媒体記憶装置を含むローカルコンピューターおよび遠隔コンピューターの記憶媒体の双方において位置し得る。またさらに、前述された命令は、一部または全部が、ハードウェア論理回路として実装されてもよく、ハードウェア論理回路は、プロセッサを含んでも、または含まなくてもよい。

10

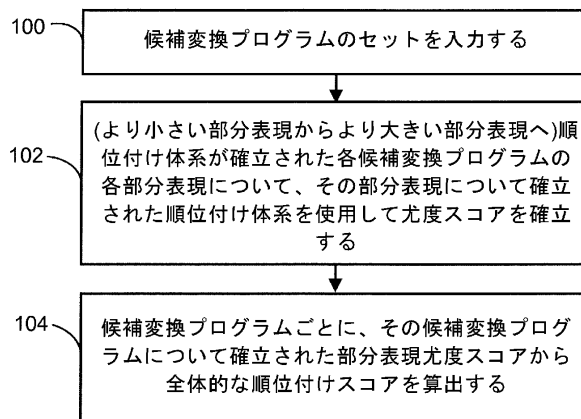
【0065】

3.0 他の実施形態

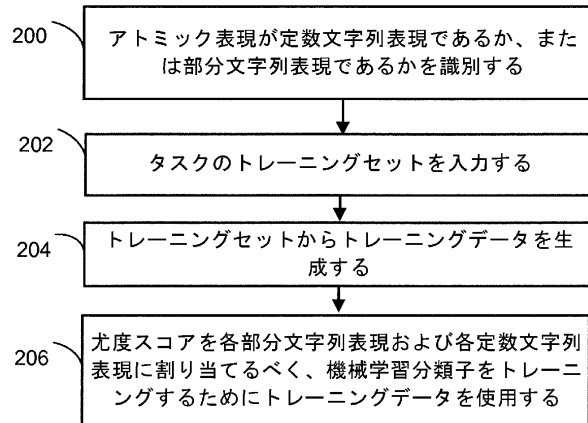
[0065] 本説明全体にわたる前述された実施形態のいずれかまたは全部は、任意の所望の組み合わせにおいて使用されて、付加的な混合実施形態を形成し得ることに留意されたい。また、主題は、構造的な特徴および/または方法論的な動作に固有の言語において説明されてきたが、添付の特許請求の範囲において定義される主題は、必ずしも上述された具体的な特徴または動作に限定されないことが理解されるべきである。むしろ、上述された具体的な特徴および動作は、請求項を実装する例示的な形式として開示される。

20

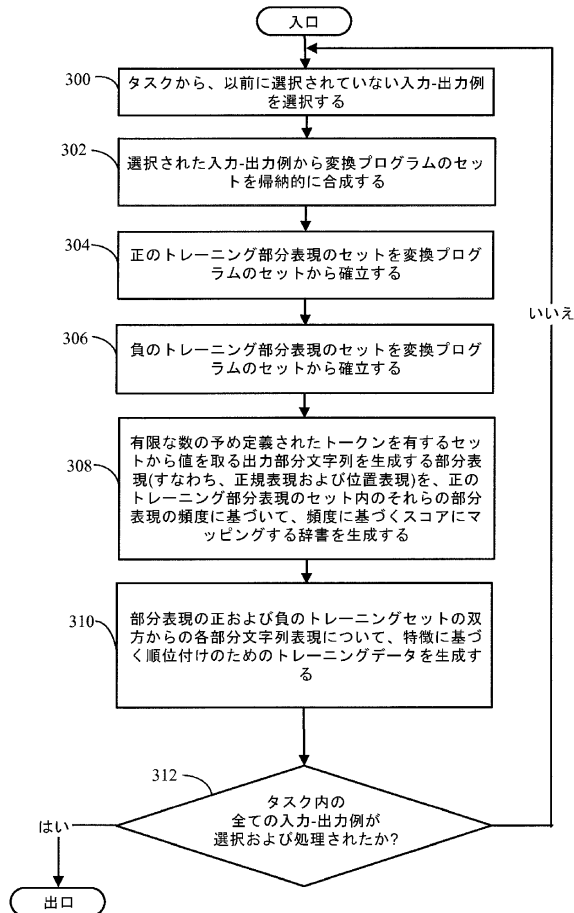
【図1】



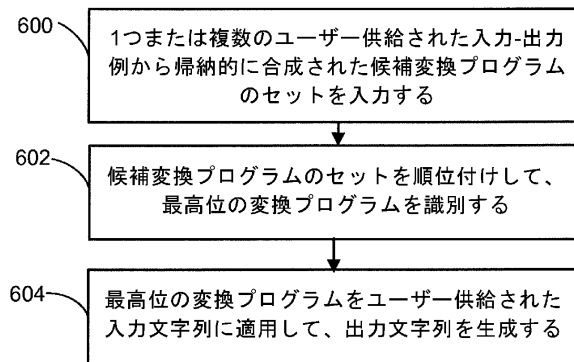
【図2】



【図 3】



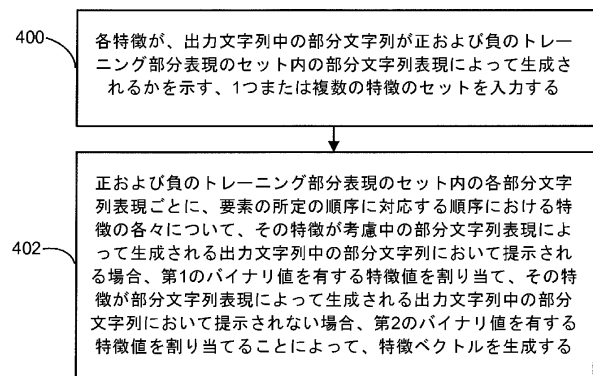
【図 6】



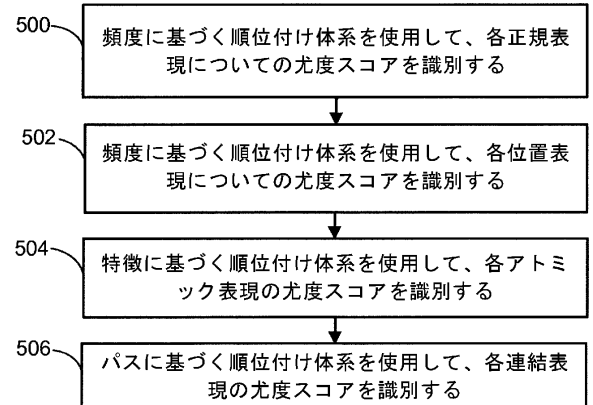
【図 7】

	入力 v_1	出力
1	Roger	Mr. Roger
2	Simon	
3	Benjamin	
4	John	

【図 4】



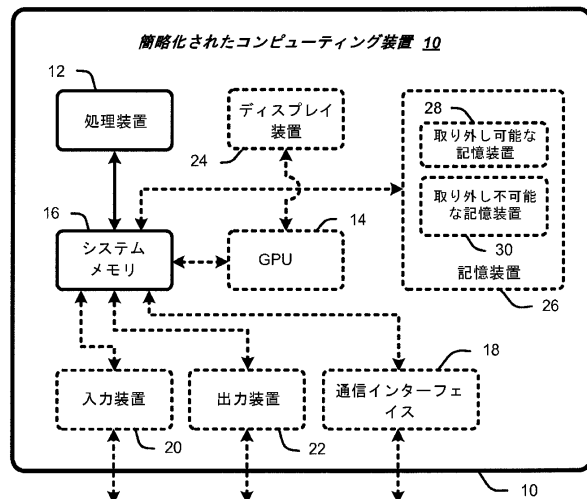
【図 5】



【図 8】

	入力 v_1	出力
1	243 Flyer Drive, Cambridge, MA 02145	Cambridge
2	512 Wright Ave, Los Angeles, CA 78911	
3	64 128th St, Seattle, WA 98102	
4	560 Hearst Ave, San Francisco, CA 94129	

【図 9】



フロントページの続き

(74)代理人 100108213

弁理士 阿部 豊隆

(74)代理人 100140431

弁理士 大石 幸雄

(72)発明者 ガルワニ, スミット

アメリカ合衆国, ワシントン州 98052-6399, レッドモンド, ワン マイクロソフト
ウェイ, マイクロソフト コーポレーション内, エルシーエー - インターナショナル パテン
ツ

(72)発明者 シン, リシャブ

アメリカ合衆国, ワシントン州 98052-6399, レッドモンド, ワン マイクロソフト
ウェイ, マイクロソフト コーポレーション内, エルシーエー - インターナショナル パテン
ツ

審査官 成瀬 博之

(56)参考文献 米国特許出願公開第2011/0302553 (US, A1)

米国特許出願公開第2012/0192051 (US, A1)

特開平07-319682 (JP, A)

(58)調査した分野(Int.Cl., DB名)

G06F 17/20 - 17/28

G06F 8/65

G06N 99/00